

GeoLinked Data

An application case/ Un caso de aplicación

Vilches Blázquez, Luis Manuel; Villazón-Terrazas, Boris; Corcho, O.; Gómez Pérez, Asunción

Resumen

La Web de los Datos enlazados, del inglés *Web of Linked Data*, supone un nuevo paradigma que pretende explotar la Web como un espacio global de información en el que la navegación se realiza a través de datos estructurados enlazados (*Linked Data*) en lugar de realizarse a través de documentos.

El término Linked Data se refiere a una forma de publicar y enlazar datos estructurados en la Web utilizando RDF (*Resource Description Framework*), un lenguaje para representar información sobre recursos propuesto por el Consorcio de la World Wide Web en el área de la Web Semántica. El valor y la utilidad de los datos enlazados es mayor cuanto más interconectados estén unos datos con otros.

La cantidad de datos enlazados publicados en la Web de los Datos ha experimentado un enorme crecimiento en los últimos años. La lista de recursos ya disponibles en *Linked Data* crece día a día. El mayor auge hasta ahora se ha producido en el contexto de la publicación de datos del sector público, siendo referencia los gobiernos del Reino Unido y de Estados Unidos. Asimismo, el fenómeno de *Linked Data* se está extendiendo a otros sectores, entre los que destacan los medios de comunicación, infraestructuras y logística, el ámbito universitario y científico. Sin embargo, hasta el momento, en el ámbito de la información geográfica los datos publicados no resultan abundantes, ya que no se han iniciado suficientes esfuerzos en esta dirección.

En este artículo se presenta el proceso seguido para el desarrollo de una aplicación que utiliza diversos conjuntos de datos relacionados con tres temas recogidos en los Anexos de la Directiva INSPIRE. Específicamente estos temas se corresponden con: unidades administrativas, hidrografía y unidades estadísticas. Esta aplicación trata de poner de manifiesto la relación existente entre la zona costera nacional e información relacionada con la población, desempleo e industria. Además, se proporcionan una guías metodológicas para la generación, publicación y explotación del *Linked Data* de las fuentes de información tratadas. Como resultado de este proceso se proporciona una importante innovación con respecto a otros procesos similares desarrollados por otras iniciativas, consecuencia de que este trabajo trata con la información geométrica de los fenómenos.

Esta transformación y publicación de la información, vinculada con diferentes temas INSPIRE, conforme a los principios de *Linked Data* añade una nueva dimensión a la Web de los Datos. De esta manera, la información geoespacial puede ser recuperada y enlazada a niveles de granularidad sin precedentes.

Abstract

In this paper we present the process that has been followed for the development of an application that makes use of several heterogeneous Spanish public datasets that are related to three themes of INSPIRE Directive, specifically Administrative Units, Hydrography, and Statistical Units. Our application aims at analysing existing relations between the Spanish coastal area and different statistical variables such as population, unemployment, dwelling, industry, and building trade. Besides providing methodological guidelines for the generation, publishing and exploitation of Linked Data from such datasets, we provide an important innovation with respect to other similar processes followed in other initiatives by dealing with the geometrical information of features.

PALABRAS CLAVE

Información geoespacial y estadística, INSPIRE, RDF, Linked Data.

KEYWORDS

Geospatial and statistical information, INSPIRE, RDF, Linked Data.

1. INTRODUCTION

The rise of the Open Data Movement has led to the Web of Data grow significantly over the last years. This Web has started to span data sources from a wide range of domains such as people, companies, music, scientific publications, etc.

Linked Data has been recently suggested as one of the best alternatives for creating these shared information spaces [9]. The notion of Linked Data refers to the recommended best practices for exposing, sharing, and connecting RDF data via dereferenceable URIs on the Semantic Web [1]. These best practices have been adopted by an increasing number of data providers, leading to the creation of a global data space containing billions of assertions - the Web of Data [9].

In the geospatial context, GeoLinked Data¹ is an open initiative whose aim is to enrich the Web of Data with Spanish geospatial data into the context of INSPIRE (*IN*frastructure for *SP*atial *IN*formation in *EU*rope) Directive² themes. This initiative has started off by publishing diverse information sources belonging to the National Geographic Institute of Spain, onwards IGN-E, and the National Statistic Institute in Spain, onwards INE. Such sources are made available as RDF knowledge bases according to the Linked Data principles.

This paper describes the results of developing an application that combines these diverse Spanish public datasets so that relationships can be inferred amongst these data. Moreover, we discuss the process followed, and propose methodological guidelines for all the activities involved within the process, which could be extrapolated to the development of similar applications.

This paper is structured as follows. We start by providing some background on the Linked Data initiative in Section 2. In Section 3 we provide a general overview of the process that we propose for publishing GeoLinked Data on the Web. The details of the activity of identification and selection of data sources are provided in Section 4. In Section 5 we present the ontology network modelling process. The generation of the RDF data is introduced in Section 6, and alignment of the datasets is detailed in Section 7. In Section 8 we describe data publication and visualization. A comparison of our approach with other geospatial Linked Data approaches is presented in Section 9. Finally, Section 10 presents some conclusions of this paper and identifies our future work.

2. AN OVERVIEW OF THE LINKED DATA INITIATIVE

Since the representation and publication of geospatial data as Linked Data is only being recently addressed, as described in the introduction, we first provide some background on the Linked Data initiative. The principles of Linked Data were first outlined by Berners-Lee in 2006 [13] using the following four guidelines:

1. Use URIs as names for things.
2. Use HTTP URIs so that people can look up those names.
3. When someone looks up a URI, provide useful information, using the standards (RDF, SPARQL).
4. Include links to other URIs, so that they can discover more things.

The Linked Data architecture described by these principles suggests that HTTP URIs may be used as resource names, whether they are electronic documents or conceptual representations (i.e. the class of Person or a particular webpage). HTTP URIs should be web resolvable, so that web clients may discover additional knowledge by following links between these web-accessible resources. This guidance has been extended by technical documents (e.g., [7][14]) that capture best practices emerging from the Linked Data community and provide recipes that can be used as a basis by Linked Data publishing systems.

Technically, in the context of Linked Data the RDF language is used to describe resources in the form of triples (subject-predicate-object), which can provide additional links (URIs in the object position of the triples), what allows connecting data between different data sources. The HTTP protocol is used to handle the interactions

¹ <http://geo.linkeddata.es/>

² The INSPIRE Directive addresses 34 spatial data themes needed for environmental applications. <http://inspire.jrc.ec.europa.eu/>

between Linked Data clients and publishers. Further details about sets of rules for publishing data on the Web are shown in [13].

3. A PROCESS FOR PUBLISHING GEOLINKED DATA ON THE WEB

¡Error! No se encuentra el origen de la referencia. provides an overview of the process that we have followed for generating and publishing GeoLinked Data from a set of data sources. This process is a generalization of the one described in ¡Error! No se encuentra el origen de la referencia., and consists of the following activities: (1) identification of the data sources, (2) generation of the ontology model, (3) generation of the RDF data, (4) publication of the RDF data, and (5) linkage of the RDF data with other existing datasets in the Web of Data.

In the following sections we will be describing each of the steps in this process, in the context of the development of an application that makes use of Spanish public datasets that are related with three themes of INSPIRE, specifically with Administrative Units, Hydrography, and Statistical Units.

The goal of this application is to analyze existing relations, in the context of the Spanish coastal area, between different statistical variables such as unemployment, population, dwelling, industry, and building trade. In this way, Open Government Data should help us to know how seasonal employment changes in these areas of Spain, where the tourism sector is very relevant for their economy.

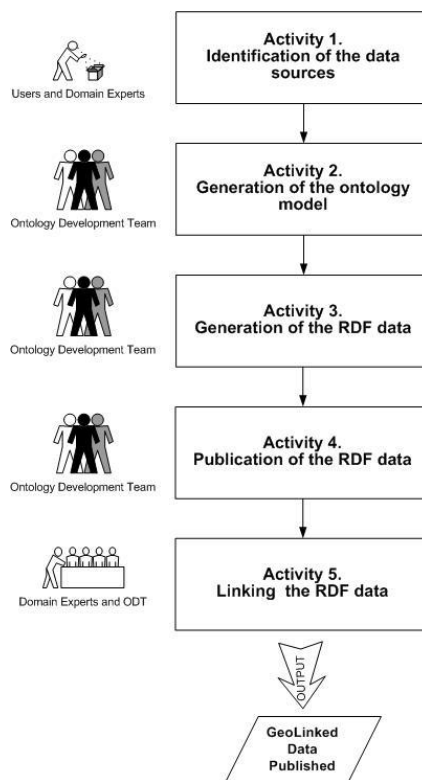


Figure 1. A process for generating and publishing GeoLinked Data

4. IDENTIFICATION AND SELECTION OF DATA SOURCES

We have searched for open government information at the two institutions that we have referred to in the introduction: INE and IGN-E. Both INE and IGN-E are providers of Spanish official statistical and geographical information, respectively. All the datasets correspond to Spain, so their content is available in Spanish or in any of the other official languages in Spain (Basque, Catalan and Galician).

Table 1. Used datasets

Data	Provider	Format
Population	INE	Excel Spreadsheet
Unemployment	INE	Excel Spreadsheet
Building Trade	INE	Excel Spreadsheet
Dwelling	INE	Excel Spreadsheet
Industry	INE	Excel Spreadsheet
Hydrography	IGN-E	Relational Database (Oracle)
Beaches	IGN-E	Relational Database (MySQL)
Adm. boundaries	IGN-E	Relational Database (MySQL)

Table 1 depicts the datasets that we have chosen for this application, together with the format in which they are available.

5. ONTOLOGY MODELLING

For the modelling of the information contained in the datasets (time, administrative boundaries, unemployment, etc.) we have created an ontology network, which is a collection of ontologies joined together through a variety of different relationships such as mapping, modularization, version, and dependency relationships [4]. This network has been developed following the NeOn methodology [5], by reusing existing ontologies and vocabularies. An overview of the GeoLinked Data ontology network is shown in Figure 2. Next, we describe briefly each one of ontologies that compose this network.

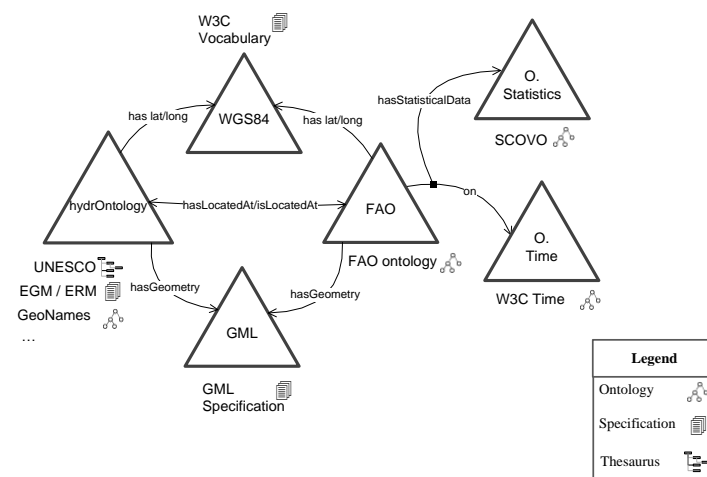


Figure 2. Ontology network of GeoLinked Data

- For representing complex statistics, we chose **Statistical Core Vocabulary (SCOVO)**³, which provides an expressive modelling framework for statistical information.
- Regarding geospatial vocabulary we chose diverse ontologies.

³ <http://vocab.deri.ie/scovo>

The **FAO geopolitical ontology**⁴. This OWL ontology includes information about continents, regions, countries and so on, in the English language. We have extended it to cover the main characteristics of the Spanish administrative division.

Regarding the hydrographical phenomena (rivers, lakes, etc.) we chose **hydrOntology**⁵, an OWL ontology build following a top-down development approach, and which attempts to cover most of the concepts of the hydrographical domain.

With respect to geometrical representation and positioning we reuse the **GML Ontology** and the **WSG84 Vocabulary**.

- Regarding the time information we chose the **Time Ontology**⁶, an ontology of temporal concepts developed into the context of World Wide Web Consortium (W3C).

Taking into account that the SCOVO and the FAO geopolitical ontologies were available in the English language, and it was important for our application to have labels in Spanish, we have used the LabelTranslator system [6] to carry out the task of ontology localization. This way we use LabelTranslator for translating components of these ontologies to Spanish.

6. GENERATION OF THE RDF DATA

We decided to use RDF as the normal form for the datasets to be published, due to the fact that we are pursuing a Linked Data approach. RDF is one of the standard languages in which information has to be made available, according to the Linked Data principles. The reason for this is that it offers several advantages, such as provision of an extensible schema, de-referenceable URIs, and as RDF links are typed, safe merging (linking) of different datasets.

Given the different formats in which the selected datasets were available, we used two different systems for the conversion of data into RDF. Next we describe some details of both of them.

The generation of RDF from spreadsheets was performed using the NOR2O [2] software library. This library performs an Extract, Transform, and Load (ETL) process of the legacy data sources, transforming these non-ontological resources (NORs) [2] into ontology instances.

The transformation of the relational database (Oracle and MySQL) content into RDF was done using the integrated framework R2O+ and ODEMapster+ [3], which is available as a NeOn Toolkit plugin⁷. This framework allows the formal specification, evaluation, verification and exploitation of semantic mappings between ontologies and relational databases.

6.1 CREATION OF RDF OF GEOMETRICAL INFORMATION

A very important aspect to be considered in this transformation of geographical information, which is also a distinctive feature of our approach, as it will be described later, is the definition of geometrical information in RDF. Next we describe our approach for transforming geometrical information into RDF.

GML and WKT. We rely on the Oracle STO UTIL package for the transformation of the geometrical data stored in the original databases into GML (*Geography Markup Language*)⁸. The generation of GML is applied to the **GEOMETRY** column, where different rows of a table have geometrical information of each feature.

For MySQL spatial databases, we work with the WKT⁹ format for extracting information of the **GEOMETRY**

⁴ <http://www.fao.org/countryprofiles/geoinfo.asp?lang=en>

⁵ <http://mayor2.dia.fi.upm.es/index.php/en/ontologies/107-hydrontology>

⁶ <http://www.w3.org/TR/owl-time/>

⁷ <http://www.neon-toolkit.org>

⁸ <http://www.opengeospatial.org/standards/gml>

⁹ *Well-Know Text* is a text markup language for representing vector geometry objects on a map, spatial reference systems of spatial objects and transformations between spatial reference systems.

example of an exploratory search interface, whose design has been investigated in some recent Human-Computer Interaction (HCI) research for supporting users who have less clear or more complex needs. The application is able to render on the map the distinct geometrical shapes of the geographical features published as RDF.

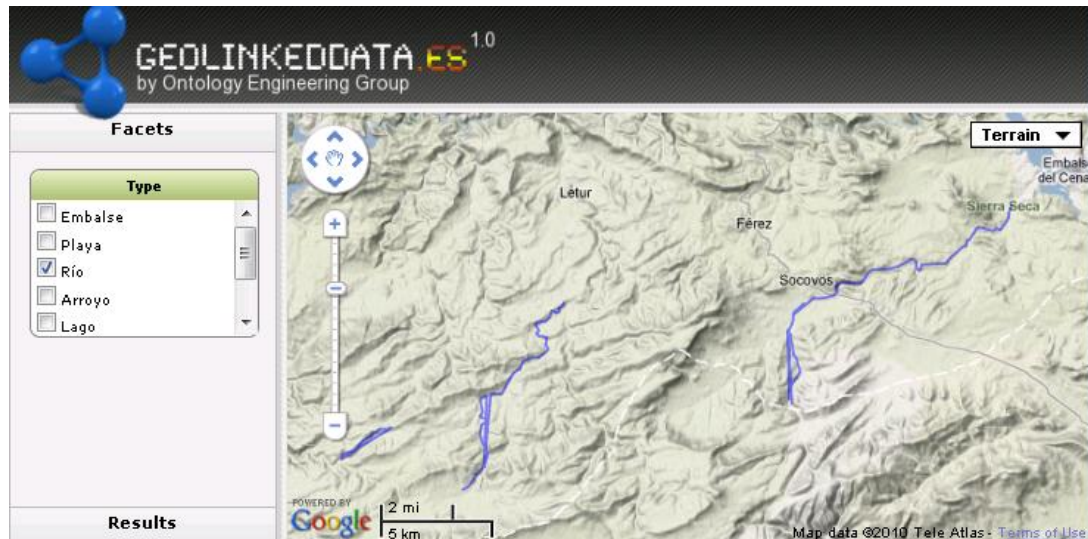


Figure 4. Visualization of rivers as StringLines

Finally, statistical data is also displayed over the map so that the user can observe and compare the relative magnitudes of the statistics, which are represented by different graphs, occurring on the distinct geographical regions, and specifically on the coastal area of the country (see Figure 5).



Figure 5. Visualization of statistical and geographical information

9. RELATED WORK

The transformation and publication of the OpenStreetMap [12] and Ordnance Survey [11] data according to the

Linked Data principles have added a new dimension to the Web of Data. Likewise, various geodata sources are starting to appear in knowledge bases of the Linked Open Data initiative such as:

- Ordnance Survey (Great Britain's national mapping agency). It provides Linked Data of the administrative and voting regions in Great Britain. This data includes the names, census code, and area of the regions of Great Britain.
- LinkedGeoData¹⁵. It provides Linked Data of the OpenStreetMap project, interlinking this data with other knowledge bases in the Linking Open Data initiative.
- GeoNames¹⁶. It integrates geographical data, such as place names, population, etc. from various sources. GeoNames gives access to users to manually edit, correct, and add new names.
- DBpedia¹⁷. It extracts structured information from Wikipedia, linking this information to other datasets, and making it available as Linked Data.

However, neither the Ordnance Survey nor the OpenStreetMap initiatives, as the main geospatial Linked Data providers to date, deal with complex geospatial information as we do in our approach. For the time being, they just manage every resource as a point, while we deal with this coordinates types and more complex geometry (LineString).

10. CONCLUSIONS AND FUTURE WORK

In this paper we have presented an application that makes use of several Spanish public datasets, specifically datasets related with three INSPIRE themes (Administrative units, Hydrography, and Statistical units). The goal of this application case is to analyze existing relations in the Spanish coastal area and different statistical variables such as unemployment, population, dwelling, industry, and building trade. Additionally, the application deals with the different geometrical information of features and establishes alignments of statistical and geometrical information. Moreover, we described the process we followed, and proposed methodological guidelines for all the activities involved.

Future work will focus on identifying and interlinking with other knowledge bases belonging to the Linking Open Data Initiative. Moreover, we will also continue publishing GeoLinked Data on the Web for other domains and providers, and improve our faceted browser. Finally, we plan to cover complex geometrical information, i.e. not only Point and LineString like data.

ACKNOWLEDGMENTS

This work has been supported by the R&D project España Virtual, funded by Centro Nacional de Información Geográfica and CDTI under the R&D programme Ingenio 2010, as well as by an R&D grant from UPM.

REFERENCES

- [1] Bizer, C., Heath, T., Idehen, K., and Berners-Lee, T. Linked data on the web (LDOW2008). In Proceeding of the 17th international conference on World Wide Web, pages 1265-1266, Beijing, China. (2008)
- [2] Villazón-Terrazas, B. Gómez-Pérez, A. and Calbimonte, J. P. NOR20: a Library for Transforming Non-Ontological Resources to Ontologies. In ESWC, volume 5554 of Lecture Notes in Computer Science. Springer. (2010)
- [3] Priyatna, F. RDF-based Access To Multiple Relational Data Sources. Master's thesis, Universidad

¹⁵ <http://linkedgeo.org/About>

¹⁶ <http://www.geonames.org/>

¹⁷ <http://dbpedia.org/>

Politécnica de Madrid. (2009)

- [4] Haase, P., Rudolph, S., Wang, Y., Brockmans, S., Palma, R., Euzenat, J., d'Aquin, M. NeOn Deliverable D1.1.1. Networked Ontology Model. Available at: <http://www.neon-project.org/>. (2006)
- [5] Suárez-Figueroa, M.C. NeOn Methodology for Building Ontology Networks: Specification, Scheduling and Reuse. PhD Thesis. Universidad Politécnica de Madrid. (2010)
- [6] Espinoza, M., Gómez-Pérez, A., and Mena, E. LabelTranslator - A Tool to Automatically Localize an Ontology. In ESWC, pages 792-796. (2008)
- [7] Bizer, C., Cyganiak, R., Heath, T. How to publish Linked Data on the Web. <http://www.wiwi-fu-berlin.de/~bizer/pub/LinkedDataTutorial/> (2007)
- [8] Oren, E., Delbru, R., and Decker, S. Extending faceted navigation for RDF data. In International Semantic Web Conference, pages 559-572. (2006)
- [9] Bizer, C., Heath, T. and Berners-Lee, T. Linked data - the story so far. International Journal on Semantic Web and Information Systems (IJSWIS). Vol. 5(3), Pages 1-22. (2009)
- [10] W3C. Publishing Open Government Data. W3C Working Draft. <http://www.w3.org/TR/gov-data/> (2009)
- [11] Goodwin, J., Dolbear, C., Hart, G. Geographical Linked Data: The Administrative Geography of Great Britain on the Semantic Web. Transactions in GIS, Volume 12 Issue s1, Pages 19 - 30 (2009)
- [12] Auer, S., Lehmann, J., Hellmann, S. LinkedGeoData - Adding a spatial Dimension to the Web of Data. In Proc. of 7th International Semantic Web Conference (ISWC), (2009)
- [13] Berners-Lee, T. Linked Data - Design Issues. W3C. <http://www.w3.org/DesignIssues/LinkedData.html> (2006)
- [14] Ayers, D., Vinkel, M. Cool URIs for the semantic web. Interest Group Note 20080331, W3C, <http://www.w3.org/TR/2008/NOTE-cooluris-20080331/> (2008)

CONTACTOS

Luis M. VILCHES-BLÁZQUEZ
lmvilches@fi.upm.es
Ontology Engineering Group
DIA - Facultad de Informática
Universidad Politécnica de
Madrid

Boris VILLAZÓN-TERRAZAS
bvillazon@fi.upm.es
Ontology Engineering Group
DIA - Facultad de Informática
Universidad Politécnica de
Madrid

Oscar CORCHO
ocorcho@fi.upm.es
Ontology Engineering Group
DIA - Facultad de Informática
Universidad Politécnica de
Madrid

Asunción GÓMEZ-PÉREZ
asun@fi.upm.es
DIA - Facultad de Informática
Universidad Politécnica de
Madrid

...

...